

# EDULISS: a small-molecule database with data-mining and pharmacophore searching capabilities

Kun-Yi Hsin, Hugh P. Morgan, Steven R. Shave, Andrew C. Hinton, Paul Taylor and Malcolm D. Walkinshaw\*

The Centre for Translational and Chemical Biology, The University of Edinburgh, King's Buildings, Edinburgh, UK

Received May 30, 2010; Revised September 3, 2010; Accepted September 17, 2010

## ABSTRACT

We present the relational database EDULISS (EDinburgh University Ligand Selection System), which stores structural, physicochemical and pharmacophoric properties of small molecules. The database comprises a collection of over 4 million commercially available compounds from 28 different suppliers. A user-friendly web-based interface for EDULISS (available at <http://eduliss.bch.ed.ac.uk/>) has been established providing a number of data-mining possibilities. For each compound a single 3D conformer is stored along with over 1600 calculated descriptor values (molecular properties). A very efficient method for unique compound recognition, especially for a large scale database, is demonstrated by making use of small subgroups of the descriptors. Many of the shape and distance descriptors are held as pre-calculated bit strings permitting fast and efficient similarity and pharmacophore searches which can be used to identify families of related compounds for biological testing. Two ligand searching applications are given to demonstrate how EDULISS can be used to extract families of molecules with selected structural and biophysical features.

## INTRODUCTION

The high throughput screening regimes of the past 20 years led by big pharma and more recently developed by screening centres through the Molecular Libraries Roadmap program are providing increasing amounts of publicly available biological information. The bioassay and compound databases in PubChem (1) contain

information on over 25 million structures and on over 60 million data points from thousands of assays. Smaller but well annotated databases like ChEMBLdb with over 500 000 entries provide information on the properties and activities of drug-like molecules and their targets (2). This explosion of data linking compounds to biological activity should provide a means for predicting new biological effects for large numbers of classes of small drug-like molecules using bioinformatic and database mining approaches (3).

In order to test such *in silico* predictions it is important to have databases of available compounds. It is only relatively recently that searchable interactive small molecule databases have become available to non-commercial research groups. One such resource is ChemDB (4), a searchable chemical database containing nearly 5 million small molecules with their stereoisomers. Interactive databases like ZINC (5) provide large and well annotated collections with some searching capacity. Such databases can contain a variety of structurally related information stored as SMILES strings, InChI or Daylight fingerprints (6). 3D coordinates may also be used as input for structure-based virtual screening (7–9) or pharmacophore searching (10). The idea of relating the activity of a molecule to the spatial distribution of a number of functional groups (11) has been widely used in QSAR (12) and structure-based studies as implemented in programs like GRID (13), LigandScout (14) and Catalyst (15).

The EDULISS database stores 3D atomic coordinates for each molecule along with over 1600 calculated molecular properties. These so called molecular descriptors provide a numerical profile for each molecule consisting of calculated values such as molecular weight, surface area and number of rotatable bonds. By using a selection of descriptors it is possible to rapidly select small related families of molecules from the database. An extension of this selection procedure provides a very efficient way of

\*To whom correspondence should be addressed. Tel: 00441316507056; Fax: 00441316507055; Email: m.walkinshaw@ed.ac.uk

identifying unique compounds. The database also stores a range of interatomic distances between various atom types for each molecule. The overall statistics of interatomic distances is used in an ultrafast shape searching algorithm (16). A specific subset of interatomic distances between all hydrogen bond donor and acceptor atoms, halogens, phosphorous and sulphur atoms provide what we call the Interatomic Pharmacophore Profile (IPP). All such distance information is stored for each molecule in pre-calculated bit-strings which provide the basis of a wide range of pharmacophore searching routines and also in the identification of similarly shaped molecules. The EDULISS database is therefore a useful tool for identifying commercially available molecules based on similarity or pharmacophore searches. It is distinguished from other web resources by having over 1600 descriptors for each compound and the ability to carry out unique 3D and 2D searches. There are also convenient links for a subset of compounds to the PubChem database allowing easy access to biological data.

## PROGRAM AND DATABASE DESCRIPTION

### Database description

Currently, EDULISS stores over 5.5 million (over 4 million unique) compounds in total, containing data from 28 different commercial and other smaller specialist compound catalogues (Supplementary Data S1). 2D and 3D coordinates for each molecule are stored with over 1600 topological, geometrical, physicochemical and toxicological descriptors per compound. In this database, over 3.9 million compounds fit the Lipinski's rule of five (17) and a total of 3.4 million fit the Oprea lead-like criteria (18): that is molecular weight  $\leq 460$ , number of rotatable bonds  $\leq 10$ , calculated Log P between  $-4$  and  $4.2$ , number of hydrogen bond acceptors  $\leq 9$ , number of hydrogen bond donors  $\leq 5$  and number of rings  $\leq 4$ . The database also contains over 520 000 compounds with molecular weight  $< 250$  Da and potentially fitting the needs of fragment-based screening (19).

The biological properties of a subset of 291 000 compounds stored in EDULISS has been retrieved from four other databases, including PubChem, BindingDB (20), ChemBank (21) and DrugBank (22), by identifying identical molecules using the Maximum Common Subgraph algorithm (23). The identity of these compounds in the external databases has been obtained and stored in the EDULISS database. A direct link between EDULISS and the external database has been implemented on the search result pages. Once a particular compound which is identical to one of the PubChem compounds has been hit by either 3D/2D similarity or molecule ID search, the link in the 'Chemical Properties' box can lead users to the appropriate PubChem web page. Certain catalogues (e.g. the National Cancer Institute) contain many compounds for which there are a lot of biological data and most hits will have links to the relevant PubChem bioassay summary page.

The EDULISS database is held in a MySQL server. The web-based interface of EDULISS uses Java Servlet

technology (see <http://java.sun.com/products/servlet/>) and JavaServer Pages (JSP, see <http://java.sun.com/products/jsp/>) to build the web pages (Figure 1). The web site utilizes Apache Tomcat as the web server and the runtime environment for Java technologies mentioned above. For the molecule drawing and visualizing, JME (<http://www.molinspiration.com/jme/>) and Jmol are utilized which are applications written in Java providing interactive features and have been included in the EDULISS web pages. On the query result page, the users can download the SDfile of hit compounds with their descriptor values. To date, this database has been used freely by the researchers from over 20 countries via its web-based interface.

### Treatment of compound structure data files, SDfiles

Regardless of the source of catalogues, the compounds used for EDULISS were entirely collected as 2D SDfile formats then converted into 3D atomic coordinates using CONCORD software. After the conversion process, the molecules were processed by DRAGON 5.4 (<http://www.taletemi.it/>) and DEREK (<http://www.lhasalimited.org>) software calculating 1664 physicochemical and potential toxicity properties for each compound.

## DATABASE APPLICATIONS

### Recognition of unique compounds

As EDULISS holds millions of compounds from various suppliers, it is useful to be able to determine the number of unique compounds in the collection. A 2D graph theory algorithm, Maximum Common Subgraph, MCS (23), has been implemented. Although the MCS is able to precisely identify isomorphous compounds, the number of pairwise comparisons increase as  $N \times (N - 1)$  where  $N$  is the number of compounds and the run time grows dramatically from 1 h to 1 day when the dataset increases from 800 to 3200 compounds (Supplementary Data S2). Thus, it is impossible to go through the whole EDULISS collection using this method.

We have developed a method to efficiently identify unique compounds by clustering according to specific descriptor values (molecular properties). Using this approach the required graphical comparisons can be considerably decreased. Preliminary studies using molecular weight and atom type were not very useful as only 6% of the compounds in EDULISS could be uniquely identified. However a number of other molecular descriptors show much better discrimination; W3D [Wiener 3D index (24)], Whete [Wiener-type index from electronegativity weighted distance matrix (25)] and Vu [a molecular size descriptor which is one of the Weighted Holistic Invariant Molecular descriptors (26)]. The combination of these three descriptors alone was sufficient to identify 3 117 625 unique compounds (out of a total of 4 011 697 unique compounds present in EDULISS). The remaining 2 million compounds were grouped using the three descriptors (W3D, Whete and Vu) into 845 193 clusters. The compounds in these clusters with identical descriptors were then compared using MCS. This procedure reduces

Descriptor Groups	Items
1 2D AUTOCORRELATIONS	96
2 3D MORSE DESCRIPTORS	160
3 ATOM CENTRED FRAGMENTS	122
4 BCUT DESCRIPTORS	64
5 CHARGE DESCRIPTORS	14
6 CONNECTIVITY INDICES	33
7 CONSTITUTIONAL DESCRIPTORS	48
8 EDGE ADJACENCY INDICES	107
9 EIGENVALUE BASED INDICES	44
10 FUNCTIONAL GROUP COUNTS	154
11 GEOMETRICAL DESCRIPTORS	74
12 GETAWAY DESCRIPTORS	197
13 INFORMATION INDICES	47
14 MOLECULAR PROPERTIES	30
15 RANDIC MOLECULAR PROFILES	41
16 RDF DESCRIPTORS	150
17 TOPOLOGICAL CHARGE INDICES	21

Submit Query

Molecular weight  
Between  ~  or  
=  OK

Determining the molecular descriptors

Suppliers: ACB-Eurochem, Ambridge, Ambrisco, Avance, ChemBridge, ChemDiv, ChemGenex, Fluka, InteBioScreen, InterChem, KeyOrganics, Lab-Chemicals (ILab), MayBridge, MayBridge\_Ro3\_Fragment, Novosyn

Num. of Compounds: 101,239, 594,832, 18,160, 149,889, 459,512, 397,434, 38,792, 382,909, 387,009, 378,207, 24,749, 180,093, 64,459, 1,237, 78,292

Chemical structure: OCC1=CC=C(C=C1)CO

Drawing the query structure

The distance between  
[Acceptor] and [Acceptor] = 5 Å Remove

The distance between  
[Acceptor] and [Acceptor] = 6 Å Remove

More criteria  
Acceptor  
Donor  
Cl  
Br  
I  
F  
S=O  
P=O

Determining the distances between specified atom types

Catalogues:	Num. of Hits:
MayBridge	61,173
Sigma-Aldrich_Family(Sigma, Aldrich, Salor, Fluka)	164,913
<b>Catalogue(s):</b>	<b>Total: 226,086</b>

You may select the following descriptor items whose values will be written into the tag block of the SD file.

<input type="checkbox"/> MW	molecular weight
<input type="checkbox"/> nCIC	number of rings
<input type="checkbox"/> REN	number of rotatable bonds
<input type="checkbox"/> nHDon	number of donor atoms for H-bonds (N and O)
<input type="checkbox"/> nHAcc	number of acceptor atoms for H-bonds (N,O,F)
<input type="checkbox"/> TPSA(Tot)	topological polar surface area using N,O,S,P polar contributions
<input type="checkbox"/> XLogP	Ghose-Crippen octanol-water partition coeff. (logP)

Download Query Again

Descriptor-based search results

Query Structure:

MW: 166.09938 | eAT: 12 | nH: 0 | nD: 2 | nP: 0 | nS: 0 | nF: 0 | nCl: 0  
RBN: 4 | nCC: 1 | nHAcc: 2 | nHDon: 2 | TPSA: 40.46 | XLogP: 0.264  
Lipinski Rule of 5: Fit | Astex Rule of 3: Not Fit | Open Lead-Like: Fit

Chemical structure: OCC1=CC=C(C=C1)CO

Calculating the molecular properties of the query structure

Chemical Properties(Descriptors)	Representation
MayBridge: AVV 00058   SPH NUMBER: SPH1-008-001 MW: 166.101   eAT: 12   nH: 0   nD: 2   nP: 0   nS: 0   nF: 0   nCl: 0 RBN: 4   nCC: 1   nHAcc: 2   nHDon: 2   TPSA: 40.46   XLogP: 0.264 Lipinski Rule of 5: Fit   Astex Rule of 3: Not Fit   Open Lead-Like: Fit Identical Compound(s): 1	<chem>OCC1=CC=C(C=C1)CO</chem>
MayBridge: AVV 00057   SPH NUMBER: SPH1-008-002 MW: 155.161   eAT: 10   nH: 0   nD: 1   nP: 0   nS: 0   nF: 0   nCl: 1 RBN: 1   nCC: 2   nHAcc: 2   nHDon: 0   TPSA: 32.38   XLogP: 1.62 Lipinski Rule of 5: Fit   Astex Rule of 3: Not Fit   Open Lead-Like: Fit Identical Compound(s): 1	<chem>OCC1=CC=C(C=C1)C</chem>
MayBridge: AVV 00054   SPH NUMBER: SPH1-008-003 MW: 278.92   eAT: 15   nH: 1   nD: 2   nP: 0   nS: 0   nF: 0   nCl: 1	<chem>OCC1=CC=C(C=C1)C2=CC=CC=C2</chem>

Showing the details of hit compounds

Chemical Properties(Descriptors)	Representation
Sigma-Aldrich_Family(Sigma, Aldrich, Salor, Fluka)   R158052   SPH NUMBER: SPH1-053-363 MW: 166.101   eAT: 12   nH: 0   nD: 2   nP: 0   nS: 0   nF: 0   nCl: 0 RBN: 4   nCC: 1   nHAcc: 2   nHDon: 2   TPSA: 40.46   XLogP: 0.26 Lipinski Rule of 5: Fit   Astex Rule of 3: Not Fit   Open Lead-Like: Fit Identical Compound(s): 0	<chem>OCC1=CC=C(C=C1)CO</chem>
Sigma-Aldrich_Family(Sigma, Aldrich, Salor, Fluka)   A7155   SPH NUMBER: SPH1-100-933 MW: 292.24   eAT: 21   nH: 0   nD: 2   nP: 0   nS: 0   nF: 0   nCl: 0 RBN: 3   nCC: 1   nHAcc: 2   nHDon: 2   TPSA: 40.46   XLogP: 4.91 Lipinski Rule of 5: Fit   Astex Rule of 3: Not Fit   Open Lead-Like: Not Fit Identical Compound(s): 0	<chem>OCC1=CC=C(C=C1)C2=CC=CC=C2</chem>
Sigma-Aldrich_Family(Sigma, Aldrich, Salor, Fluka)   S310719   SPH NUMBER: SPH1-096-166 MW: 194.131   eAT: 14   nH: 0   nD: 2   nP: 0   nS: 0   nF: 0   nCl: 0	<chem>OCC1=CC=C(C=C1)C</chem>

Structure-based similarity search results

Jmol

Supplier	Original ID	SPH_NUMBER
1 InterBioScreen	STOCK38-23979	SPH1-210-032
2 ChemBridge	6938468	SPH1-257-807
3 Sigma-Aldrich_Family(Sigma, Aldrich, Salor, Fluka)	L352314	SPH1-131-349

Listing the identical compounds

Search options

Output

Figure 1. The EDULISS web-based interface provides four search options, including descriptor-based searches, structure-based similarity searches, IPP searches and search by molecule ID. On the query result pages, the users can download the SDfile of hit compounds with their descriptor values.

the number of required pair-wise comparisons using MCS down to 6495096 which can be carried out in 20 h.

### Similarity searching

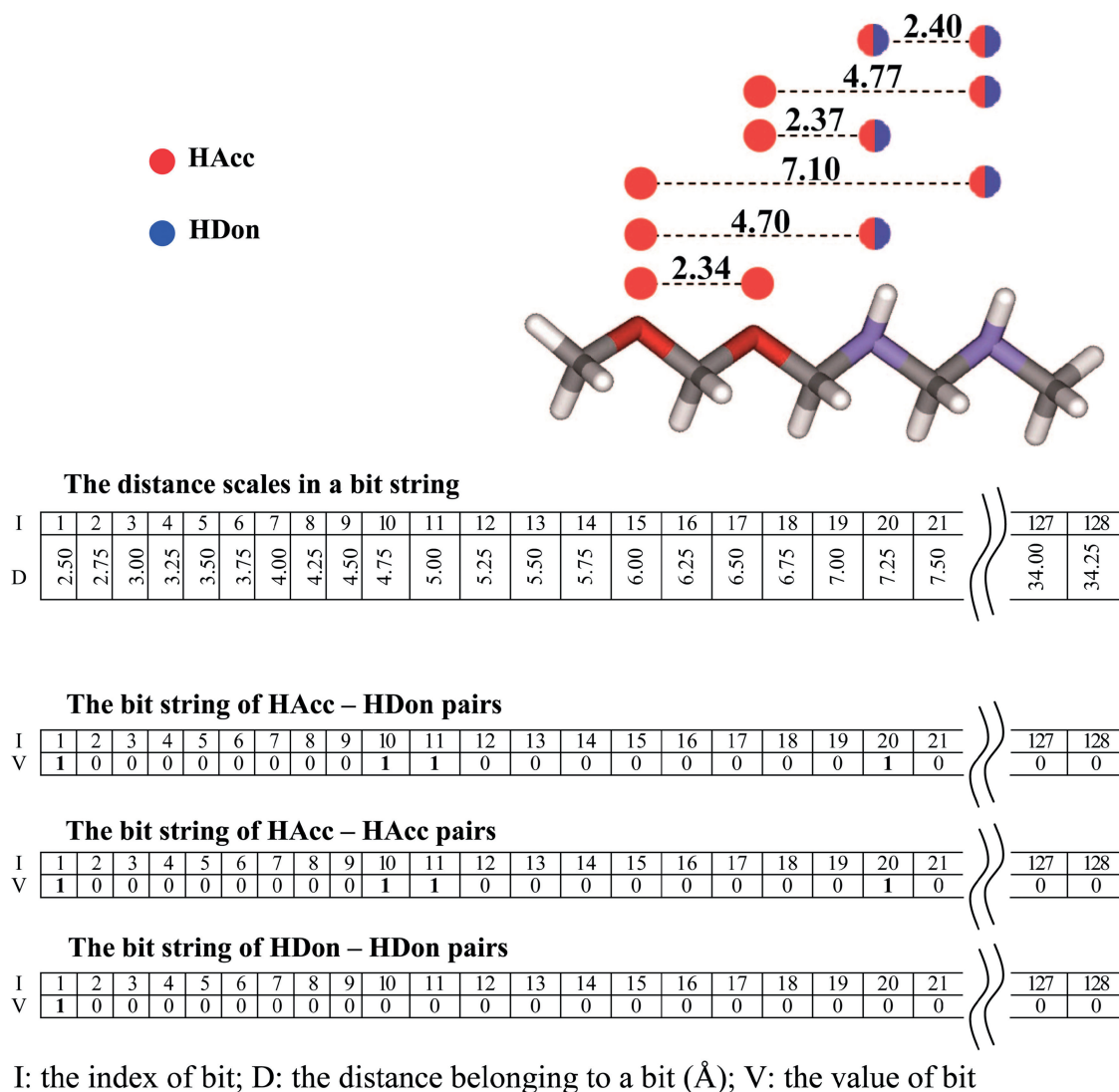
EDULISS stores more than 1600 molecular descriptors for each compound and users can select a series of descriptor items as a query to identify a subset of molecules which will share common properties. Molecular descriptors are primarily organized into 20 groups according to their attributes, so that the users can conveniently choose and set preferred values for the query. For example it is a simple matter to extract from the Sigma-Aldrich catalogue the 164913 out of 199492 compounds that fit the Lipinski rule of five and the 142660 that comply with the Oprea lead-like criteria.

EDULISS also provides geometrical similarity searches based on a 3D similarity measurement called Ultra Fast Shape Recognition with Atom Types (UFSRAT). UFSRAT uses pre-generated geometric descriptors for

molecules within EDULISS to discriminate between both the overall geometric, hydrophobic and electrostatic shape of molecules.

### Pharmacophore searching using EDULISS

The IPP for each molecule in the database consists of interatomic distances calculated between 8 different atom classes; namely hydrogen bond donor atoms (HDon), hydrogen bond acceptor atoms (HAcc), halogens (fluorine, chlorine, bromine and iodine), sulphur and phosphorus atoms. This gives rise to 15 possible types of interatomic distance for each molecule. Distances are stored in strings 128 bits long as Boolean values (1, true; 0, false). The first bit represents a distance less than or equal to 2.50 Å, the next bit is 0.25 Å longer (i.e. >2.50 and ≤2.75 Å) and so forth until the last bit which represents any distance >34.00 Å. Thus, there are 15 bit strings for each molecule representing the 15 types of possible distance pairs. Figure 2 illustrates the



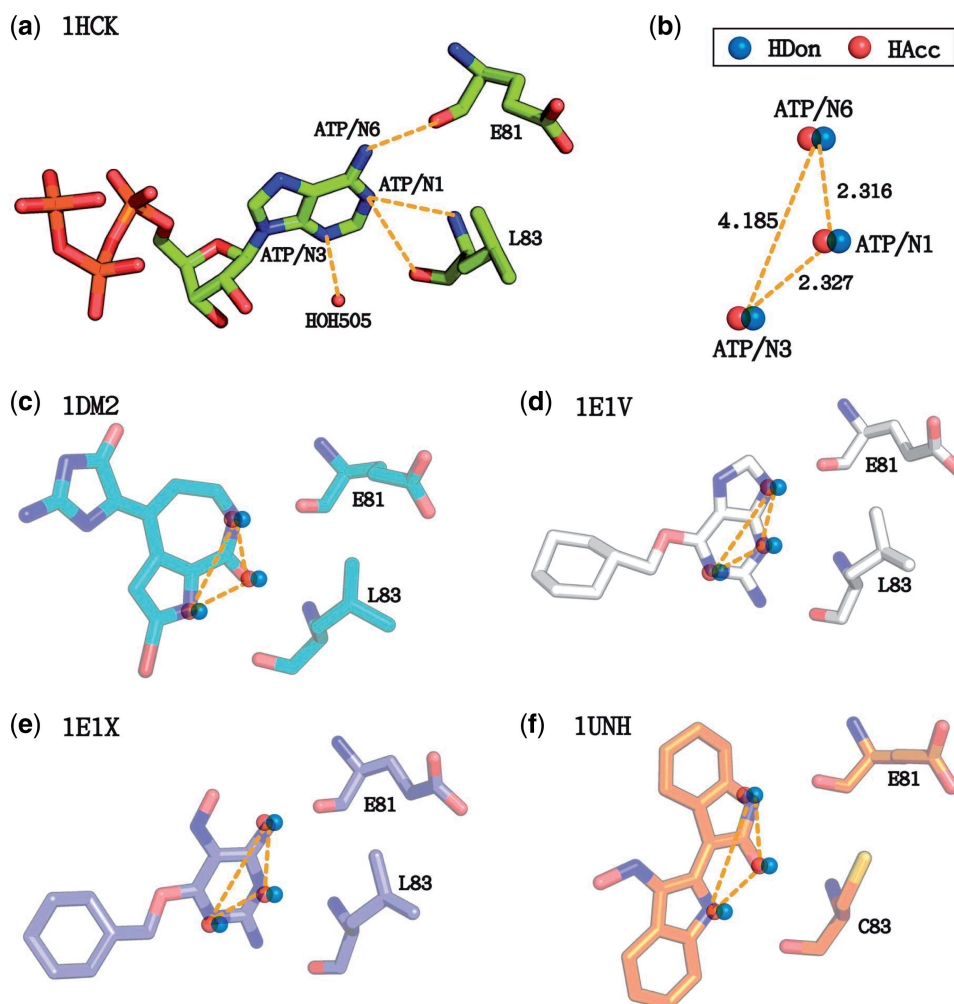
**Figure 2.** Examples of the bit string composition of a virtual compound. Hydrogen bond acceptor (HAcc) atoms are coloured red and hydrogen bond donor (HDon) atoms are coloured blue.

composition of three bit strings showing distances between HAcc and HDon.

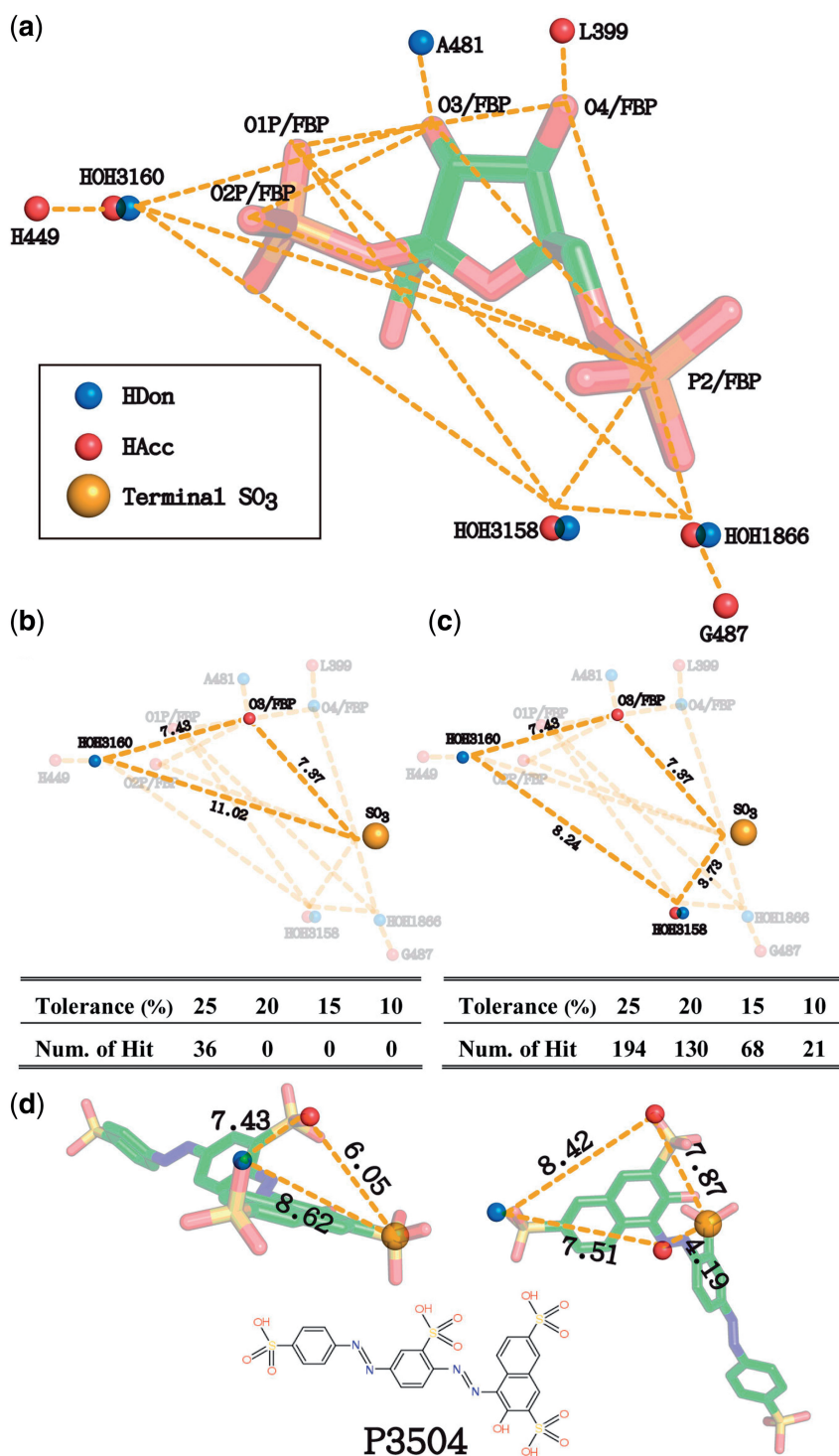
This facility enables compounds to be identified that have a specific geometric arrangement of atoms ('pharmacophore') as defined by pair-wise distances of hydrogen bond donors, acceptors, halogens, phosphorous or sulphur atoms (where the searches are restricted to S and P atoms that form double bonds to oxygen). For pharmacophore searching, a bit string is generated for the user-defined query distances which are then compared to that of each compound in the database. If a specific true bit in the query matches, the distance criterion is met. A user can perform a multi-distance query in a single search. Apart from very efficient storage, bit strings also provide a very fast searching method as the necessary Boolean operations can be carried out very quickly. Users can specify the query by defining preferred distances between selected atom types using the web-based interface. The results are then displayed and hits may be downloaded.

*Case study 1: identifying cyclin dependent kinases inhibitors.* As a test for the pharmacophore searching routine we used the eight available structures of CDK complexes stored in the PDBbind database (27) with PDB codes 1AQ1, 1DI8, 1DM2, 1E1V, 1E1X, 1FVV, 1UNH and 2A4L. Three nitrogen atoms of the adenine ring of ATP were used as a template to search for ATP-analogues (Figure 3). Applying the three interatomic distance criteria as shown in Figure 3b, four out of the eight ligands (1DM2, 1E1V, 1E1X and 1UNH) were identified as illustrated in Figure 3c–f, (The other four ligands were not recognized as they do not have the adenine-like pharmacophore)

*Case study 2: identifying pyruvate kinase inhibitors.* The glycolytic enzyme pyruvate kinase (PYK) is a drug target against trypanosomatid infection (28). Fructose-2,6-bisphosphate (F-2,6-BP) acts as an allosteric activator (29). We are interested in identifying analogue molecules which interfere with allosteric regulation. Figure 4a schematically shows selected interatomic



**Figure 3.** Schematic diagrams illustrating the CDK inhibitor pharmacophore search. (a) the key interaction of ATP-binding pocket in the CDK2–ATP complex (PDB id: 1HCK). (b) interatomic distances between the three atoms selected from the ATP adenine ring. The hydrogen bond donors (HDon) and acceptors (HAcc) are coloured blue and red and labelled by residue name and ID. (c–f) show the interesting interactions of the four hit CDK proteins whose ligand-bound interactions are similar to the pharmacophore search model.



**Figure 4.** Schematic diagrams of the relevant interactions of F-2,6-BP (FBP) in the effector site of PYK. (a) shows the interatomic distances between five atoms in F-2,6-BP and three water molecules, selected for pharmacophore models. (b and c) Examples of pharmacophore search models with hydrogen bond donors (HDon) and acceptors (HAcc) coloured blue and red. Sulfonate (SO<sub>3</sub>) is coloured yellow. (d) illustrates the hit compound P3504 in different orientations to show that both pharmacophore profiles fit.

contacts and distances between five atoms in F-2,6-BP and three water molecules selected for pharmacophore searches. Two example search motifs are shown in Figure 4b and c. A series of tolerances have been given for each interatomic distance from 10 to 25%. The

numbers of hit compounds in a range of tolerances are tabulated in Figure 4. We selected eight compounds (Sigma-Aldrich ID: N9002, L0144, P3504, 201332, 244813, D5021, H2516 and 86170) for further experimental assay based on visual inspection of the docked pose

and on calculated solubility. Of the eight selected compounds, five significantly affected the PYK enzyme kinetics. Figure 4d shows that both pharmacophore search models match atoms of the hit compound P3504 which showed 33% inhibition of enzyme activity. A complex of P3504 with *Leishmania mexicana* PYK (LmPYK) has been crystallized and solved at a resolution of 2.7 Å (H. P. Morgan personal communication) showing the molecule binds at the effector site.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

The authors would like to thank the staff at Centre of Translational and Chemical Biology.

## FUNDING

The Wellcome Trust and the Scottish University Life Sciences Alliance for the use of Edinburgh Protein Production Facility. K.-Y.H. who did most of the work was supported by an Edinburgh University departmental scholarship. This work was not directly supported by any granting agencies or by commercial companies.

*Conflict of interest statement.* None declared.

## REFERENCES

- Geer,L.Y., Marchler-Bauer,A., Geer,R.C., Han,L., He,J., He,S., Liu,C., Shi,W. and Bryant,S.H. (2010) The NCBI BioSystems database. *Nucleic Acids Res.*, **38**, D492–D496.
- EMBL-EBI. (2010) Open access drug discovery database launches with half a million compounds: [www.ebi.ac.uk/chembl/db/index.php](http://www.ebi.ac.uk/chembl/db/index.php). *EMBL-EBI Press Release*.
- Miller,M.A. (2002) Chemical database techniques in drug discovery. *Nat. Rev. Drug Discov.*, **1**, 220–227.
- Chen,J.H., Linstead,E., Swamidass,S.J., Wang,D. and Baldi,P. (2007) ChemDB update—full-text search and virtual chemical space. *Bioinformatics (Oxford, England)*, **23**, 2348–2351.
- Irwin,J.J. and Shoichet,B.K. (2005) ZINC: A free database of commercially available compounds for virtual screening. *J. Chem. Inf. Model.*, **45**, 177–182.
- Butina,D. (1999) Unsupervised data base clustering based on daylight's fingerprint and Tanimoto similarity: A fast and automated way to cluster small and large data sets. *J. Chem. Inf. Comput. Sci.*, **39**, 747–750.
- Anderson,A.C. (2003) The process of structure-based drug design. *Chem. Biol.*, **10**, 787–797.
- Ghosh,S., Nie,A., An,J. and Huang,Z. (2006) Structure-based virtual screening of chemical libraries for drug discovery. *Curr. Opin. Chem. Biol.*, **10**, 194–202.
- Lyne,P.D. (2002) Structure-based virtual screening: an overview. *Drug Discov. Today*, **7**, 1047–1055.
- Fang,X. and Wang,S. (2002) A web-based 3D-database pharmacophore searching tool for drug discovery. *J. Chem. Inf. Comput. Sci.*, **42**, 192–198.
- McGregor,M.J. and Muskal,S.M. (2000) Pharmacophore fingerprinting. 2. Application to primary library design. *J. Chem. Inf. Comput. Sci.*, **40**, 117–125.
- Deanda,F. and Stewart,E.L. (2004) Application of the PharmPrint methodology to two protein kinases. *J. Chem. Inf. Comput. Sci.*, **44**, 1803–1809.
- Kastenholz,M.A., Pastor,M., Cruciani,G., Haaksma,E.E. and Fox,T. (2000) GRID/CPCA: a new computational tool to design selective ligands. *J. Med. Chem.*, **43**, 3033–3044.
- Wolber,G. and Langer,T. (2005) LigandScout: 3-D pharmacophores derived from protein-bound ligands and their use as virtual screening filters. *J. Chem. Inf. Model.*, **45**, 160–169.
- Patel,Y., Gillet,V.J., Bravi,G. and Leach,A.R. (2002) A comparison of the pharmacophore identification programs: Catalyst, DISCO and GASP. *J. Comput. Aided Mol. Des.*, **16**, 653–681.
- Ballester,P.J. and Richards,W.G. (2007) Ultrafast shape recognition to search compound databases for similar molecular shapes. *J. Comput. Chem.*, **28**, 1711–1723.
- Lipinski,C.A., Lombardo,F., Dominy,B.W. and Feeney,P.J. (1997) Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Adv. Drug Deliv. Rev.*, **23**, 3–25.
- Hann,M.M. and Oprea,T.I. (2004) Pursuing the leadlikeness concept in pharmaceutical research. *Curr. Opin. Chem. Biol.*, **8**, 255–263.
- Hartshorn,M.J., Murray,C.W., Cleasby,A., Frederickson,M., Tickle,I.J. and Jhoti,H. (2005) Fragment-based lead discovery using X-ray crystallography. *J. Med. Chem.*, **48**, 403–413.
- Liu,T., Lin,Y., Wen,X., Jorissen,R.N. and Gilson,M.K. (2007) BindingDB: a web-accessible database of experimentally determined protein-ligand binding affinities. *Nucleic Acids Res.*, **35**, D198–D201.
- Seiler,K.P., George,G.A., Happ,M.P., Bodycombe,N.E., Carrinski,H.A., Norton,S., Brudz,S., Sullivan,J.P., Muhlich,J. and Serrano,M. (2008) ChemBank: a small-molecule screening and cheminformatics resource database. *Nucleic Acids Res.*, **36**, D351–D359.
- Wishart,D.S., Knox,C., Guo,A.C., Cheng,D., Shrivastava,S., Tzur,D., Gautam,B. and Hassanali,M. (2007) DrugBank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Res.*, **36**(Suppl. 1), D901–D906.
- Raymond,J.W. and Willett,P. (2002) Maximum common subgraph isomorphism algorithms for the matching of chemical structures. *J. Comput. Aided Mol. Des.*, **16**, 521–533.
- Mihalic,Z. and Veljan,D. (1992) The distance matrix in chemistry. *J. Math. Chem.*, **11**, 223–258.
- Ivanciuc,O., Ivanciuc,T. and Balaban,A.T. (1998) Design of topological indices. Part 10. Parameters based on electronegativity and covalent radius for the computation of molecular graph descriptors for heteroatom-containing molecules. *J. Chem. Inf. Comput. Sci.*, **38**, 395–401.
- Todeschini,R., Lasagni,M. and Marengo,E. (1994) New molecular descriptors for 2D and 3D structures. *Theory. J. Chemom.*, **8**, 263–272.
- Wang,R., Fang,X., Lu,Y. and Wang,S. (2004) The PDBbind database: collection of binding affinities for protein-ligand complexes with known three-dimensional structures. *J. Med. Chem.*, **47**, 2977–2980.
- Mesecar,A.D. and Nowak,T. (1997) Metal-ion-mediated allosteric triggering of yeast pyruvate kinase. 2. A multidimensional thermodynamic linked-function analysis. *Biochemistry*, **36**, 6803–6813.
- Morgan,H.P., McNaie,I.W., Hsin,K.Y., Michels,P.A.M., Fothergill-Gilmore,L.A. and Walkinshaw,M.D. (2010) An improved strategy for the crystallization of *Leishmania mexicana* pyruvate kinase. *Acta Crystallogr. Sect. F Struct. Biol. Cryst. Commun.*, **66**, 215–218.